

(11) Publication number: **09311802 A**

**(43) Date of publication of application: 02.12.97**

(51) Int. Cl. **G06F 12/00**  
**G06F 12/00**  
**G06F 11/34**  
**G06F 15/00**  
**G06F 17/30**

(21) Application number: 08149784

**(22) Date of filing: 22.05.96**

(71) Applicant: **MATSUSHITA ELECTRIC IND CO LTD**

(72) Inventor: **UENO TAKESHI  
NOGUCHI YOSHIHIRO  
SATO MITSUHIRO  
ISHIKAWA MIKITO**

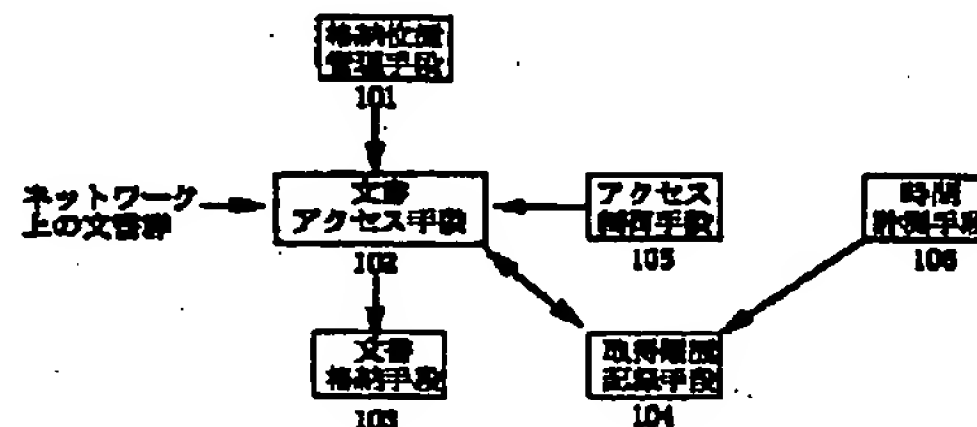
**(54) DOCUMENT GATHERING DEVICE**

**(57) Abstract:**

**PROBLEM TO BE SOLVED:** To provide a document gathering device for automatically and efficiently gathering always newest document information at the time of gathering document in a network.

**SOLUTION:** By the start from an access control means 105, a document access means 102 gathers a desired document from a document group on a network by using a storing position managing means 101 managing the correspondence of the document position on the network and the document name, a document storing means stores it, and a time measuring means 106 records time at that time in an obtained history recording means 104. At the time of gathering document after then, only documents updated after the time recorded in the recording means 104 are gathered.

**COPYRIGHT: (C)1997,JPO**



(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平9-311802

(43) 公開日 平成9年(1997)12月2日

| (51) Int.Cl. <sup>8</sup>     | 識別記号  | 庁内整理番号 | F I           | 技術表示箇所  |
|-------------------------------|-------|--------|---------------|---------|
| G 0 6 F 12/00                 | 5 1 7 |        | G 0 6 F 12/00 | 5 1 7   |
|                               | 5 4 5 |        |               | 5 4 5 M |
| 11/34                         |       |        | 11/34         | C       |
| 15/00                         | 3 1 0 |        | 15/00         | 3 1 0 U |
| 17/30                         |       |        | 15/401        | 3 4 0 B |
| 審査請求 未請求 請求項の数 6 F D (全 10 頁) |       |        |               |         |

(21) 出願番号 特願平8-149784

(22) 出願日 平成8年(1996)5月22日

(71) 出願人 000005821

松下電器産業株式会社

大阪府門真市大字門真1006番地

(72) 発明者 上野 剛

大阪府門真市大字門真1006番地 松下電器  
産業株式会社内

(72) 発明者 野口 喜伴

大阪府門真市大字門真1006番地 松下電器  
産業株式会社内

(72) 発明者 佐藤 光弘

大阪府門真市大字門真1006番地 松下電器  
産業株式会社内

(74) 代理人 弁理士 役 昌明 (外2名)

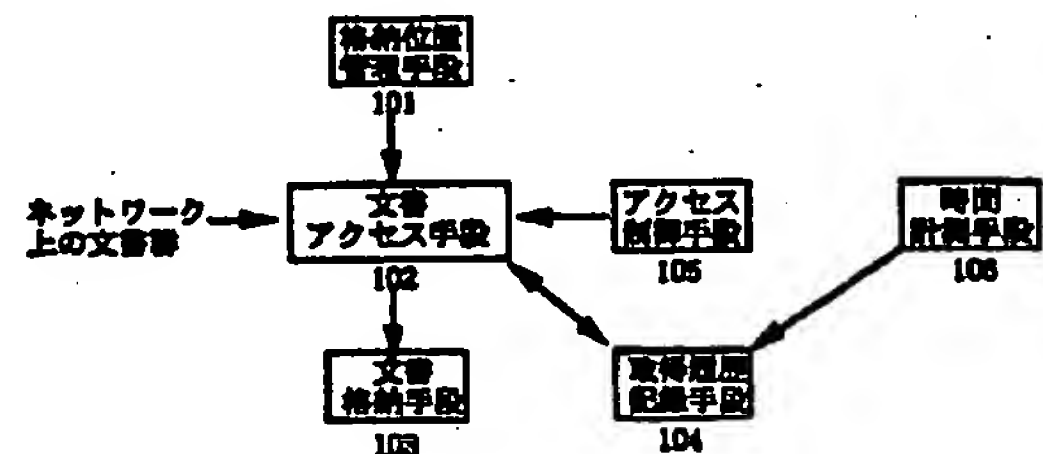
最終頁に続く

(54) 【発明の名称】 文書収集装置

(57) 【要約】

【課題】 ネットワーク上にある文書を収集する時に、常に最新の文書情報を自動的にかつ効率良く収集するための文書収集装置を提供すること。

【解決手段】 アクセス制御手段105からの起動により、ネットワーク上の文書位置と文書名の対応を管理する格納位置管理手段101を用いて文書アクセス手段102がネットワーク上の文書群から目的の文書を収集し、文書格納手段103に格納し、時間計測手段106により、その時の時刻を取得履歴記録手段106に記録し、その後の文書収集で前記取得履歴記録手段104に記録された時刻以降に更新された文書のみを収集する。



## 【特許請求の範囲】

【請求項1】 ネットワーク上に分散格納されそれぞれ別個に更新される文書群のデータを該ネットワークに接続された特定の計算機内に自動的に収集する文書収集装置において、文書収集をする文書アクセス手段と、それを駆動するアクセス制御手段をもうけ、前回収集した時点以降に更新された文書のみを収集することを特徴とする文書収集装置。

【請求項2】 文書が格納されているネットワーク上の位置と文書の対応を管理する格納位置管理手段と、ネットワーク上の文書の内容を計算機内に読み込む文書アクセス手段と、読み込んだ文書の内容を記憶する文書格納手段と、時間を計測する時間計測手段と、前記文書を読み込んだ日時を前記時間計測手段からの時刻で記録する取得履歴記録手段と、前記文書を最後に取得した取得日時とネットワーク上での前記文書の更新日時を比較し前回収集した時点以降に更新された文書のみを収集するように前記文書アクセス手段を制御するアクセス制御手段から構成されることを特徴とする文書収集装置。

【請求項3】 文書とともに得られた更新日時を更新履歴として記録する更新履歴記録手段を更に備え、前記更新履歴記録手段に記録された文書毎の複数の更新日時の記録により、前記アクセス制御手段が前記文書の平均的な更新周期を計算し前記文書を前記更新周期毎かつ予測される更新日時の直後に取得するように前記時間計測手段を利用して前記文書アクセス手段を制御することを特徴とする請求項2記載の文書収集装置。

【請求項4】 前記更新履歴記録手段に記録された文書毎の複数の更新日時の記録により更新パターンを抽出する更新パターン抽出手段と、その更新パターンを記録する更新パターン記録手段とを更に備え、前記アクセス制御手段が前記文書を前記更新パターンに合わせかつ予測される更新日時の直後に取得するように前記文書アクセス手段を制御することを特徴とする請求項3記載の文書収集装置。

【請求項5】 文書のアクセス時間と前記文書のアクセスに関する平均データ伝送率を記録するデータ伝送率記録手段を更に備え、前記文書アクセス制御手段が平均データ伝送率の大きい文書を優先して取得するように前記文書アクセス手段を制御することにより効率的な収集を可能にしたことを特徴とする請求項2記載の文書収集装置。

【請求項6】 前記データ伝送率記録手段に記録された前記文書のアクセスに関する平均データ伝送率に基づき、前記文書アクセス制御手段が平均データ伝送率の異なる複数の文書へのアクセスを組み合わせたアクセスプランを作成し、前記文書アクセス手段が複数の文書を並行にアクセスすることにより、ネットワーク資源から得られる可能な限りのデータ伝送率をより均一かつ効率的に利用できるようにしたことを特徴とする請求項5記載の

文書収集装置。

## 【発明の詳細な説明】

【0001】

【発明の属する技術分野】 本発明は、ネットワーク上に配置された文書群のデータを該ネットワークに接続された特定の計算機内に自動的に収集する文書収集装置に関し、特に文書の取得履歴を記録し、その取得履歴に基づき最新文書の取得を自動的に行なう文書収集装置に関するものである。

【0002】

【従来の技術】 従来、ネットワーク上の文書群のデータを自動的に収集する方法として、図11に示すネットワークロボットと呼ばれる文書収集手段がある。このネットワークロボットは、ネットワーク上の文書群の格納位置を与えて置くと、その文書群のデータを自動的に収集してくる機構である。このネットワークロボットの構成は図12のようになる。

【0003】 また、特開平6-301577号公報のように、情報源に固有の知識を持ち、その情報源に文書収集装置を派遣し最新の情報を得るものがある。

【0004】

【発明が解決しようとする課題】 しかしながら、前記従来のネットワークロボットでは、指定された文書を指定された順序に従ってすべて収集するため、ネットワーク上での文書の更新のパターンによっては、古いデータを収集してしまう、前回収集時以降更新されていないために収集する必要のない文書を収集してしまう、文書へのアクセスに関わるデータ伝送率などの資源を有効に利用できないなどの欠点があった。

【0005】 また、情報源に文書収集装置を派遣するものは、派遣された情報源で情報の更新の有無をチェックするため、複数の情報源を考慮した場合の収集の最適化や資源の有効利用が困難であり、文書収集の効率化に限界があった。

【0006】 本発明は、前記従来の課題を解決するもので、ネットワーク上にある文書を収集する際に、常に最新の文書情報を自動的にかつ効率良く収集するための文書収集装置を提供することを目的とする。

【0007】

【課題を解決するための手段】 この課題を解決するため、本発明は、文書収集をする文書アクセス手段と、それを駆動するアクセス制御手段をもうけ、文書の取得履歴を記録し、その取得履歴に基づき、最新文書の取得を自動的に行なうようにする。

【0008】

【発明の実施の形態】 本発明の請求項1記載の発明は、ネットワーク上に分散格納された文書群を自動的に収集する文書収集装置において、文書収集をする文書アクセス手段と、それを駆動するアクセス制御手段をもうけ、前回収集した時点以降に更新された文書のみを収集する

ことを特徴とする文書収集装置であり、最新の文書を効率的に収集保持できるという作用を有する。

【0009】本発明の請求項2記載の発明は、文書が格納されているネットワーク上の位置と文書の対応を管理する格納位置管理手段と、ネットワーク上の文書の内容を読み込む文書アクセス手段と、読み込んだ文書の内容を記憶する文書格納手段と、時間を計測する時間計測手段と、前記文書を読み込んだ日時を前記時間計測手段からの時刻で記録する取得履歴記録手段と、前記文書を最後に取得した取得日時とネットワーク上での前記文書の更新日時を比較して前回収集した時点以降に更新された文書のみを収集するために前記文書アクセス手段を制御するアクセス制御手段から構成されることを特徴とする文書収集装置であり、請求項1記載の発明と同様な作用を有する。

【0010】本発明の請求項3記載の発明は、文書とともに得られた更新日時を更新履歴として記録する更新履歴記録手段を更に備え、前記更新履歴記録手段に記録された文書毎の複数の更新日時の記録により、前記アクセス制御手段が前記文書の平均的な更新周期を計算し前記文書を前記更新周期毎かつ予測される更新日時の直後に取得するように前記時間計測手段を利用して前記文書アクセス手段を制御することを特徴とする請求項2記載の文書収集装置としたものであり、適切な更新時刻を予想して文書収集することができるため、一度に文書収集する場合に比べ負荷の分散をすることができるという作用と、常に最新の文書を収集保持することができるという作用を有する。

【0011】本発明の請求項4記載の発明は、前記更新履歴記録手段に記録された文書毎の複数の更新日時の記録により更新パターンを抽出する更新パターン抽出手段と、その更新パターンを記録する更新パターン記憶手段とを更に備え、前記文書を前記更新パターンに合わせかつ予測される更新日時の直後に取得するように前記文書アクセス手段を制御することを特徴とする請求項3記載の文書収集装置としたものであり、更新パターンを抽出し、その更新パターン毎に文書を収集することで、更新パターンにあった時だけ文書収集し、平均的に最新の文書収集ができるという作用を有する。

【0012】本発明の請求項5に記載の発明は、文書のアクセス時間と前記文書のアクセスに関する平均データ伝送率を記録するデータ伝送率記録手段を更に備え、前記文書アクセス制御手段が平均データ伝送率の大きい文書を優先して取得するように前記文書アクセス手段を制御することにより効率的な収集を可能にしたことを特徴とする請求項2記載の文書収集装置としたものであり、平均データ伝送率の小さいいくつかの文書を収集しないことなどにより、効率的な文書収集が可能であるという作用を有する。

【0013】本発明の請求項6に記載の発明は、前記デ

ータ伝送率記録手段に記録された前記文書のアクセスに関する平均データ伝送率に基づき、前記文書アクセス制御手段が平均データ伝送率の異なる複数の文書へのアクセスを組み合わせたアクセスプランを作成し、前記文書アクセス手段が複数の文書を並行にアクセスすることにより、ネットワーク資源から得られる可能な限りのデータ伝送率をより均一かつ効率的に利用できるようにしたことを特徴とする請求項5記載の文書収集装置としたものであり、前記ネットワーク資源から得られる可能な限りのデータ伝送率をより均一かつ効率的に利用できるとい

う作用を有する。

【0014】以下、本発明の実施の形態について、図1から図10を用いて説明する。

【0015】(第1の実施の形態) 本発明の第1の実施の形態について、図1を参照しながら説明する。図1は本発明の第1の実施の形態における文書収集装置の構成を示す概念図である。図1において、文書収集装置は、文書が格納されているネットワーク上の位置と文書の対応を管理する格納位置管理手段101と、ネットワーク上の文書の内容を読み込む文書アクセス手段102と、読み込んだ文書の内容を記憶する文書格納手段103と、前記文書を読み込んだ日時を後記する時間計測手段106からの時刻で記録する取得履歴記録手段104と、前記文書を最後に取得した取得日時とネットワーク上での前記文書の更新日時を比較して前回収集した時点以降に更新された文書のみを収集するために前記文書アクセス手段102を制御するアクセス制御手段105と、時間を計測する時間計測手段106とから構成されている。

【0016】以上のように構成された文書収集装置について、以下その動作を説明する。図2は文書収集装置の動作手順を示す。

【0017】まず、ステップ202において、アクセス制御手段105が文書アクセス手段102を起動する。ところで、格納位置管理手段101には文書アクセス手段102がアクセスするネットワークでの文書の位置と文書名が格納されている。

【0018】例えば、文書名、文書位置の順に、  
"文書A http://a/b/cl.html"  
のような記述が複数格納されているとする。

【0019】また、履歴記録手段104は、前回に文書を取得した文書名とその取得日時が、例えば、  
"文書A Got 1996-03-18-10:55:30"  
のように記録されている。

【0020】ステップ203では、格納位置管理手段101の文書名、文書位置の組を全てチェックしたか判定する。まだ、何もチェックしていないので、ステップ204に進む。

【0021】ステップ204において、文書アクセス手段102は、格納位置管理手段101から文書名、文書位置を一つ得る。また、履歴記録手段104から文書名にマッチす

10

20

30

40

50



る取得時間を得て、それらを使用してネットワークから目的の文書を得る。例えば

" 文書A http://a/b/cl.html" と " 文書A Got 1996-03-18-10:55:30"

から文書格納手段103に、前回取得した日時以降に更新された文書のみを得る。

【0022】また、履歴記録手段104に時間計測手段106から現在の時刻を得て、

" 文書A Got 1996-04-18-10:55:30"

のようにその内容を更新する。

【0023】この動作を、ステップ203において、格納位置管理手段101内の文書名、文書位置の組を全てチェックし終わるまで繰り返す。

【0024】以上のように、第1の実施の形態によれば、最新の文書収集を自動的に行なってこれを保持することができるという効果を有する。

【0025】(第2の実施の形態) 本発明の第2の実施の形態について、図3を参照しながら説明する。図3は本発明の第2の実施の形態における文書収集装置の構成を示す概念図である。図3において、文書収集装置は、文書が格納されているネットワーク上の位置と文書の対応を管理する格納位置管理手段301と、ネットワーク上の文書の内容を読み込む文書アクセス手段302と、読み込んだ文書の内容を記憶する文書格納手段303と、前記文書を読み込んだ日時を後記する時間計測手段306からの時刻で記録する取得履歴記録手段304と、前記文書を最後に取得した取得日時とネットワーク上での前記文書の更新日時を比較して前回収集した時点以降に更新された文書のみを収集するために前記文書アクセス手段302を制御するアクセス制御手段305と、時間を計測する時間計測手段306と、文書が更新された日時を記録する更新履歴記録手段307とから構成されている。

【0026】以上のように構成された文書収集装置について、以下その動作を説明する。図4は文書収集装置の動作手順を示す。

【0027】まず、更新履歴記録手段307は、文書名とその文書が更新された日時の履歴を、例えば

" 文書A Modified 1996-03-01-10:01:00  
1996-04-01-10:01:00  
1996-05-01-10:01:00"

のように記録しているものとする。

【0028】そこで、ステップ402において、アクセス制御手段305は、更新履歴記録手段307内の更新データから、平均的な更新周期を計算し、最後に更新された日時にその平均周期を加えた時刻(予想更新時刻)を計算する。

【0029】次にステップ403において、時間計測手段306から現在の時刻を得て、予想更新時刻を超過したか判定する。

【0030】そして、ステップ404において、予想更新

時刻を超過したら、その文書を文書アクセス手段302を起動することでアクセスする。例えば、この例では、文書Aは、平均更新周期は約1カ月なので、最後の  
" Modified 1996-05-01-10:02:00"

に1カ月を加えた

" Modified 1996-06-01-10:02:00"

の時刻(予想更新時刻)を超過したら、目的の文書Aに文書アクセス手段302を起動してアクセスする。

【0031】ここで、格納位置管理手段301には、文書アクセス手段302がアクセスするネットワークでの文書の位置と文書名を格納されているものとする。例えば、文書名、文書位置の順に、

" 文書A http://a/b/cl.html"

のような記述が複数格納されているとする。

【0032】また、取得履歴記録手段304には、前回に文書を取得した文書名とその取得日時が、

" 文書A Got 1996-03-18-10:55:30"

のように記録されているものとする。

【0033】ステップ405において、文書アクセス手段302が格納位置管理手段301からアクセス制御手段305で指定された文書名と文書位置を一つ得る。取得履歴記録手段304から文書名にマッチする取得時間を得て、それらを使用してネットワークから目的の文書を得る。例えば、

" 文書A http://a/b/cl.html" と " Got 1996-03-18-10:55:30"

から文書格納手段303に前回取得した日時以降に更新された文書のみを得る。

【0034】また、取得履歴記録手段304に時間計測手段306から現在の時刻を得て

" 文書A Got 1996-06-01-10:03:30"

のようにその内容を更新する。また、更新履歴記録手段307に取得文書の更新日時を追加する。そして、ステップ402から繰り返す。

【0035】以上のように、第2の実施の形態によれば、常に最新の文書収集を自動的に行なってこれを保持することができるという効果を有する。

【0036】(第3の実施の形態) 本発明の第3の実施の形態について、図5を参照しながら説明する。図5は本発明の第3の実施の形態における文書収集装置の構成を示す概念図である。図5において、文書収集装置は、文書が格納されているネットワーク上の位置と文書の対応を管理する格納位置管理手段501と、ネットワーク上の文書の内容を読み込む文書アクセス手段502と、読み込んだ文書の内容を記憶する文書格納手段503と、前記文書を読み込んだ日時を時間計測手段506からの時刻で記録する取得履歴記録手段504と、前記文書を最後に取得した取得日時とネットワーク上での該文書の更新日時を比較して前回収集した時点以降に更新された文書のみを収集するために文書アクセス手段502を制御するア

セス制御手段505と、時間を計測する時間計測手段506と、文書が更新された日時を記録する更新履歴記録手段507と、更新履歴記録手段507から各文書の更新パターンを抽出する更新パターン抽出手段508と、更新パターン抽出手段508が抽出した更新パターンを記録する更新パターン記録手段509とから構成されている。

【0037】以上のように構成された文書収集装置について、以下その動作を説明する。図6は文書収集装置の動作手順を示す。

【0038】まず、更新履歴記録手段507は、文書名とその文書が更新された日時の履歴を、例えば

" 文書A Modified 1996-03-01-10:01:00  
1996-04-01-10:01:00  
1996-05-01-10:01:00"

のように記録しているものとする。

【0039】そしてステップ602において、更新パターン抽出手段508は、更新履歴記録手段507内の文書履歴から、各文書の平均的な更新周期を計算し、同じ更新周期に対応する各文書を更新パターンとして、更新パターン記録手段509に記録する。

【0040】ところで、更新パターン記録手段509は、更新パターンと文書名を、

" 1年毎:文書B,文書C" と " 1ヶ月:文書A"

のように記録しているものとする。これは、文書B、文書Cは1年毎に、文書Aは1ヶ月毎に更新していることを表す。

【0041】ステップ603において、アクセス制御手段505は、更新パターン記録手段509から更新周期が短いものから順に更新周期とこの更新周期に属する文書名を得る。さらに更新履歴記録手段507から一致する文書名で、最後に更新された日時に、この更新周期を加えた日時(予想更新時刻)を計算する。これを全ての更新時刻について行ない、得られた予想更新時刻を時刻の早いものから順にソートする。

【0042】ステップ604において、時間計測手段506から得た現在時刻は、ソートされた予想更新時刻にあるか判定する。

【0043】ステップ605において、予想更新時刻を経過したものがあれば、その文書を文書アクセス手段502を起動することでアクセスする。例えば、この例では、文書Aは平均更新周期は約1ヶ月なので、最後の

" Modified 1996-05-01-10:02:00"

に1ヶ月を加えた

" Modified 1996-06-01-10:02:00"

の時刻(予想更新時刻)を経過したら、目的の文書Aに文書アクセス手段502を起動してアクセスする。

【0044】ここで、格納位置管理手段501には、文書アクセス手段502がアクセスするネットワークでの文書の位置と文書名を格納されているものとする。例えば、文書名、文書位置の順に、

" 文書A http://a/b/cl.html"

のような記述が複数格納されているとする。

【0045】また、取得履歴記録手段504には、前回に文書を取得した文書名とその取得日時が、

" 文書A Got 1996-03-18-10:55:30"

のように記録されているものとする。

【0046】ステップ606において、文書アクセス手段502が格納位置管理手段501からアクセス制御手段505で指定された文書名と一致する文書位置を得る。また、取得履歴記録手段504から文書名にマッチする取得時間を得て、それらを使用してネットワークから目的の文書を得る。例えば、

" 文書A http://a/b/cl.html" と " 文書A Got 1996-03-18-10:55:30"

から文書格納手段503に前回取得した日時以降に更新された文書のみを得る。また、取得履歴記録手段504に時間計測手段506から現在時刻を得て、

" 文書A Got 1996-06-01-10:03:30"

のようにその内容を更新する。また、更新履歴記録手段507に取得文書の更新日時を追加する。そして、ステップ603から繰り返す。

【0047】以上のように、第3の実施の形態によれば、更新パターンにあった時だけ文書収集を行ない、平均的に最新の文書収集を効率よく自動的に行なってこれを保持することができるという効果を有する。

【0048】(第4の実施の形態) 本発明の第4の実施の形態について、図7を参照しながら説明する。図7は本発明の第4の実施の形態における文書収集装置の構成を示す概念図である。図7において、文書収集装置は、文書が格納されているネットワーク上の位置と文書の対応を管理する格納位置管理手段701と、ネットワーク上の文書の内容を読み込む文書アクセス手段702と、読み込んだ文書の内容を記憶する文書格納手段703と、前記文書を読み込んだ日時を後記する時間計測手段706からの時刻で記録する取得履歴記録手段704と、後記するデータ伝送率記録手段707に記録された前記文書のアクセス時間と平均データ転送率に基づき文書アクセス手段702を制御するアクセス制御手段705と、時間を計測する時間計測手段706と、前記文書を読み込んだ際のアクセス時間と平均データ転送率を記録するデータ伝送率記録手段707とから構成されている。

【0049】以上のように構成された文書収集装置について、以下その動作を説明する。図8は文書収集装置の動作手順を示す。ここで、データ伝送率記録手段707には、各文書への過去のアクセスに基づき、文書のアクセス時間と該文書のアクセスに関する平均データ伝送率が記録されているものとする。

【0050】ステップ802において、まず、アクセス制御手段705は、データ伝送率記録手段707に記録されている平均データ伝送率が大きな文書から順にアクセスする

10

20

30

40

50



ように、文書アクセス手段702を起動する。この時、時間計測手段706よりアクセス開始時刻を得る。

【0051】ステップ803において、文書アクセス手段702は格納位置管理手段701から、指定された文書名に一致する文書位置を得る。また、取得履歴記録手段704から、この文書の前回得た日時を得る。そして、この文書をネットワーク上の文書群から指定された文書位置で、前回得た日時以降に更新された文書のみを得て、文書格納手段703に格納する。この時、時間計測手段706から現在の時刻を得て、取得履歴記録手段704に記録する。

【0052】ステップ804において、アクセス制御手段705は、時間計測手段706からアクセス終了時刻を得て、アクセス開始時刻からの経過時間を計算する。また、文書格納手段703から取得文書のサイズを得る。これらの文書サイズと経過時間から、データ伝送率を計算し、データ伝送率記録手段707に記録する。

【0053】ステップ805において、データ伝送率記録手段707内の全ての文書をチェックしたか判定する。ステップ805でチェックしていなければ、ステップ802から繰り返す。ステップ805でチェックしていれば、ステップ806に進み終了する。

【0054】以上のように、第4の実施の形態によれば、アクセス制御手段は、平均データ伝送率の順に各文書を読み込むよう文書アクセス手段を制御するので、これにより、収集すべき全文書中の大部分をより早く収集できる。また、平均データ伝送率の小さいいくつかの文書を収集しないことにより、効率的な文書収集が可能である。

【0055】（第5の実施の形態）本発明の第5の実施の形態について、図9を用いて説明する。本発明の第5の実施の形態における文書収集装置の構成は図7に示される前記第4の実施の形態の文書収集装置の構成と変わらない。したがって、文書収集装置の動作について説明する。図9は文書収集装置の動作手順を示す。ここでデータ伝送率記録手段707には、各文書への過去のアクセスに基づき、文書のアクセス時間と該文書のアクセスに関する平均データ伝送率が記録されているものとする。

【0056】ステップ902において、まず、アクセス制御手段705は、データ伝送率記録手段707に記録されている平均データ伝送率に基づき、アクセスプランを作成する。アクセスプランとは平均データ伝送率が異なる文書へのアクセスを組み合わせることで、ネットワーク資源から得られる可能な限りのデータ伝送率をより均一かつ効率的に利用するものである。アクセスプランの作成方法には、いわゆる組み合わせ最適化理論に基づく各種の方法がありうるが、すべての場合に最良解を出す方法はない。また、実際に文書を収集すると、最初に予測したアクセス時間と平均データ伝送率とは一般に差異を生ずるため、段階的にアクセスプランを修正して行く必要

がある。

【0057】図10を用いてアクセスプラン作成方法の一例を説明する。このアクセスプランの目的は、与えられたデータ伝送率という資源を最大限利用して、全文書を最短の時間で収集することにある。

【0058】まず、予想されるアクセス時間が最長の文書にアクセスする。同時に、残ったデータ伝送率を超えて最も近い平均データ伝送率を持つ文書にアクセスする。それがなければ、平均データ伝送率が最大の文書にアクセスし、再び残ったデータ伝送率に対して同様の文書を探す。利用できるデータ伝送率を超えると、文書の実際のデータ伝送率は少し低下するが、ネットワークが提供するデータ伝送率という資源は最大限利用される。

【0059】一つの文書の収集が終了した時点で、まだ収集していない文書中から、再び予想されるアクセス時間が最長の文書にアクセスし、残ったデータ伝送率を超えて最も近い平均データ伝送率を持つ文書にアクセスする。すべての文書を収集するまで、同様の手順を続ける。

【0060】図10の例では文書Aから文書Fまでの6つの文書を収集している。文書Bと文書Eの実際のデータ伝送率は、利用できるデータ伝送率の上限に達しているために少し低下している。

【0061】このように作成されたアクセスプランに基づき、ステップ903において、アクセス制御手段705は、1つ以上の文書アクセス手段702を文書名を指定して起動する。すなわち並列にアクセス手段702を起動し、文書収集を並列に行なう。この時、時間計測手段706から現在の時刻を得てアクセス開始時間を得る。

【0062】ステップ904において、文書アクセス手段702は格納位置管理手段701から、指定された文書名に一致する文書位置を得る。また、取得履歴記録手段704から、この文書を前回得た日時を得る。そして、この文書をネットワーク上の文書群から指定された文書位置で、前回得た日時以降に更新された文書のみを得て、文書格納手段703に格納する。この時、時間計測手段706から現在の時刻を得て、取得履歴記録手段704に記録する。

【0063】ステップ905において、アクセス制御手段705は、時間計測手段706からアクセス終了時刻を得て、アクセス開始時刻からの経過時間を計算する。また、文書格納手段703から取得文書のサイズを得る。これらの文書サイズと経過時間から、データ伝送率を計算し、データ伝送率記録手段707に記録する。

【0064】ステップ906において、アクセスプランを全て実行したか判定する。ステップ906でチェックしていなければ、ステップ903から繰り返す。ステップ906でチェックしていれば、ステップ907に進み終了する。

【0065】以上のように、第5の実施の形態によれば、アクセス制御手段は、アクセスプランに基づいて各文書を読み込んで行くように文書アクセス手段を制御す

10

20

30

40

50

るので、これにより、収集すべき全文書中をより早く収集でき、効率的な文書収集が可能である。

【0066】

【発明の効果】 以上のように、本発明は、最新の文書収集を自動的に行なうことができる。また、ネットワーク上の文書の次の更新時刻を予測して文書収集をすることができ、常に最新の文書収集を効率的に自動的に行なうことができる。

【0067】 また、ネットワーク上の文書の更新パターンを得て更新パターン毎に文書収集することができるので、さらに効率的に文書収集ができる。

【0068】 また、平均データ伝送率の大きい文書を優先して取得することができるので、平均データ伝送率の小さいいくつかの文書を収集しないことなどにより、効率的な文書収集ができる。

【0069】 また、文書アクセス手段が複数の文書を並行にアクセスすることができる場合に、ネットワーク資源から得られる可能な限りのデータ伝送率を利用することで、効率的な文書収集ができる。

【図面の簡単な説明】

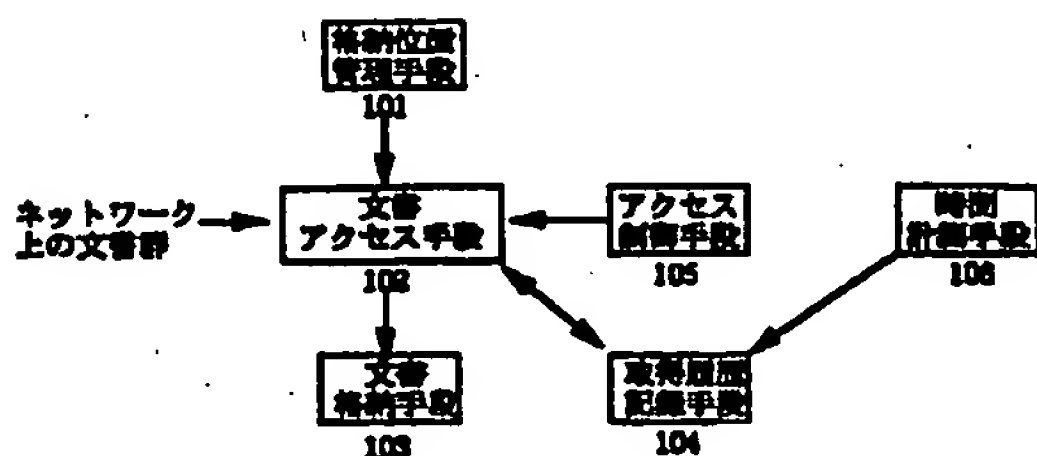
【図1】 第1の実施の形態における文書収集装置の構成を示すブロック図、

【図2】 第1の実施の形態における文書収集の手順を示すフロー図、

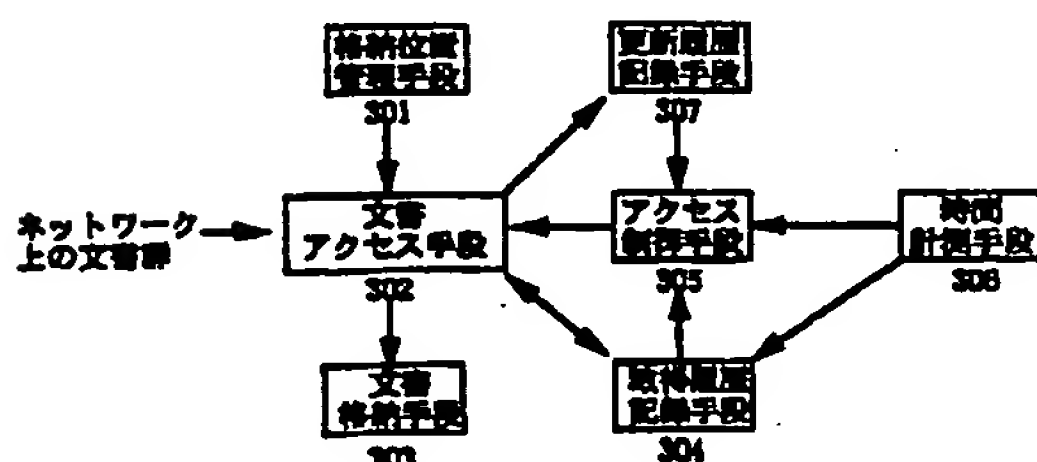
【図3】 第2の実施の形態における文書収集装置の構成を示すブロック図、

【図4】 第2の実施の形態における文書収集の手順を示すフロー図、

【図1】



【図3】



【図5】 第3の実施の形態における文書収集装置の構成を示すブロック図、

【図6】 第3の実施の形態における文書収集の手順を示すフロー図、

【図7】 第4及び第5の実施の形態における文書収集装置の構成を示すブロック図、

【図8】 第4の実施の形態における文書収集の手順を示すフロー図、

【図9】 第5の実施の形態における文書収集の手順を示すフロー図、

【図10】 アクセスプラン作成方法の一例を示す概念図、

【図11】 ネットワークロボットの使用環境を示す概念図、

【図12】 従来システムの構成を示すブロック図である。

【符号の説明】

101、301、501、701、1201 格納位置管理手段

102、302、502、702、1202 文書アクセス手段

20 103、303、503、703、1203 文書格納手段

104、304、504、704 取得履歴記録手段

105、305、505、705 アクセス制御手段

106、306、506、706 時間計測手段

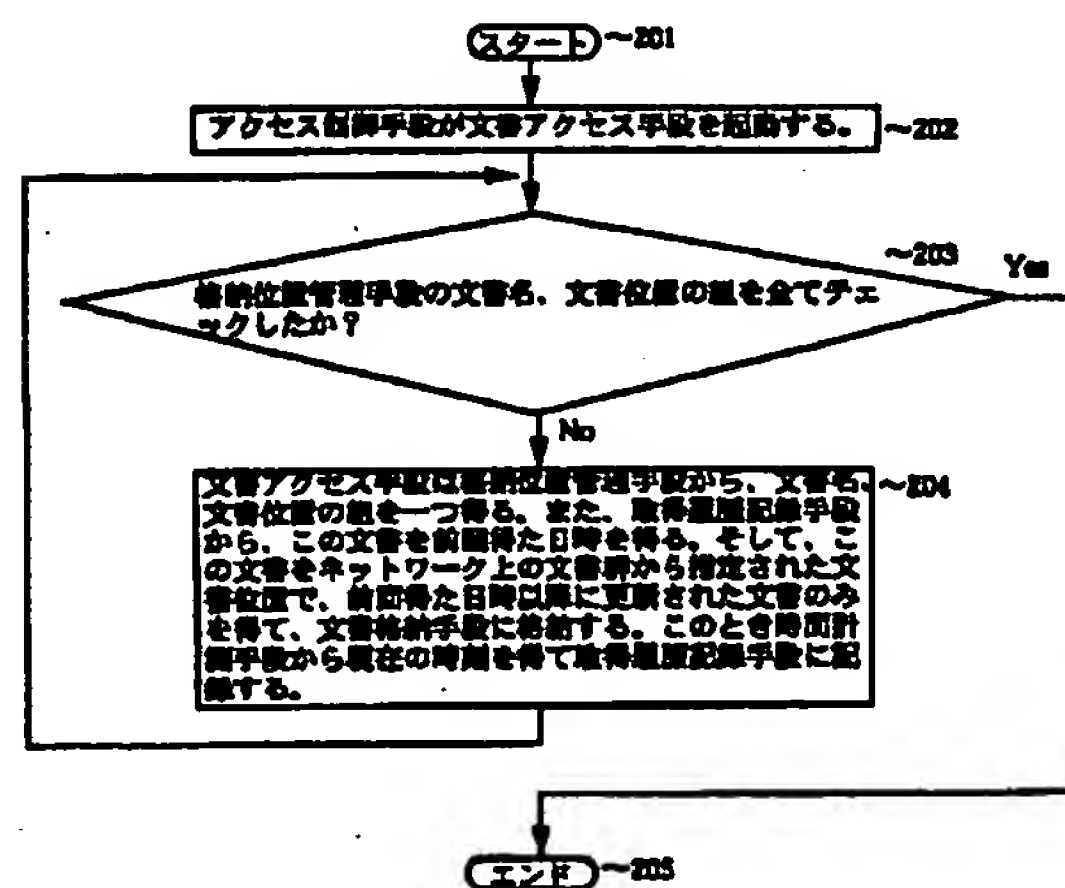
307、507 更新履歴記録手段

508 更新パターン記録手段

509 更新パターン抽出手段

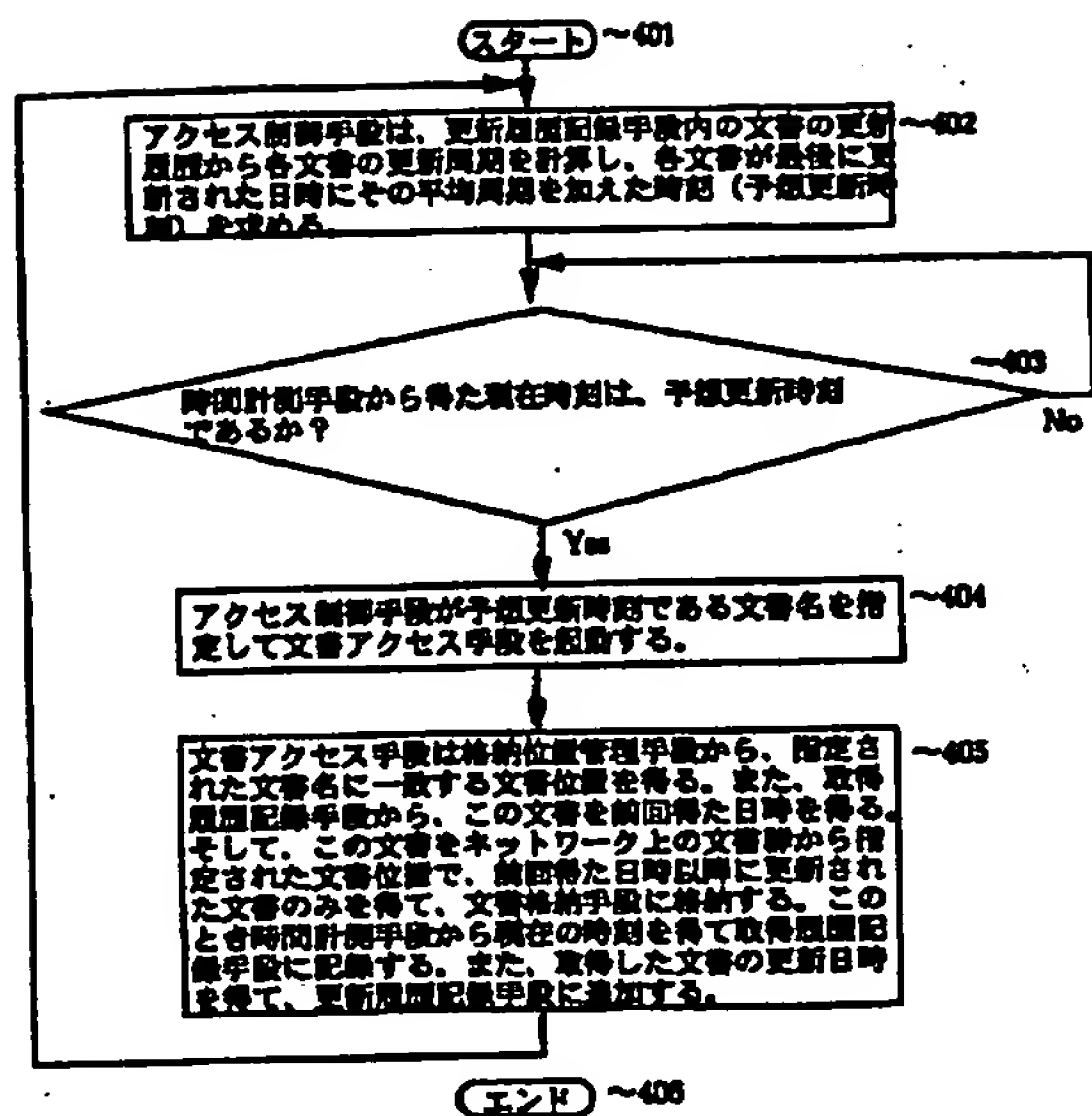
707 データ伝送率記録手段

【図2】

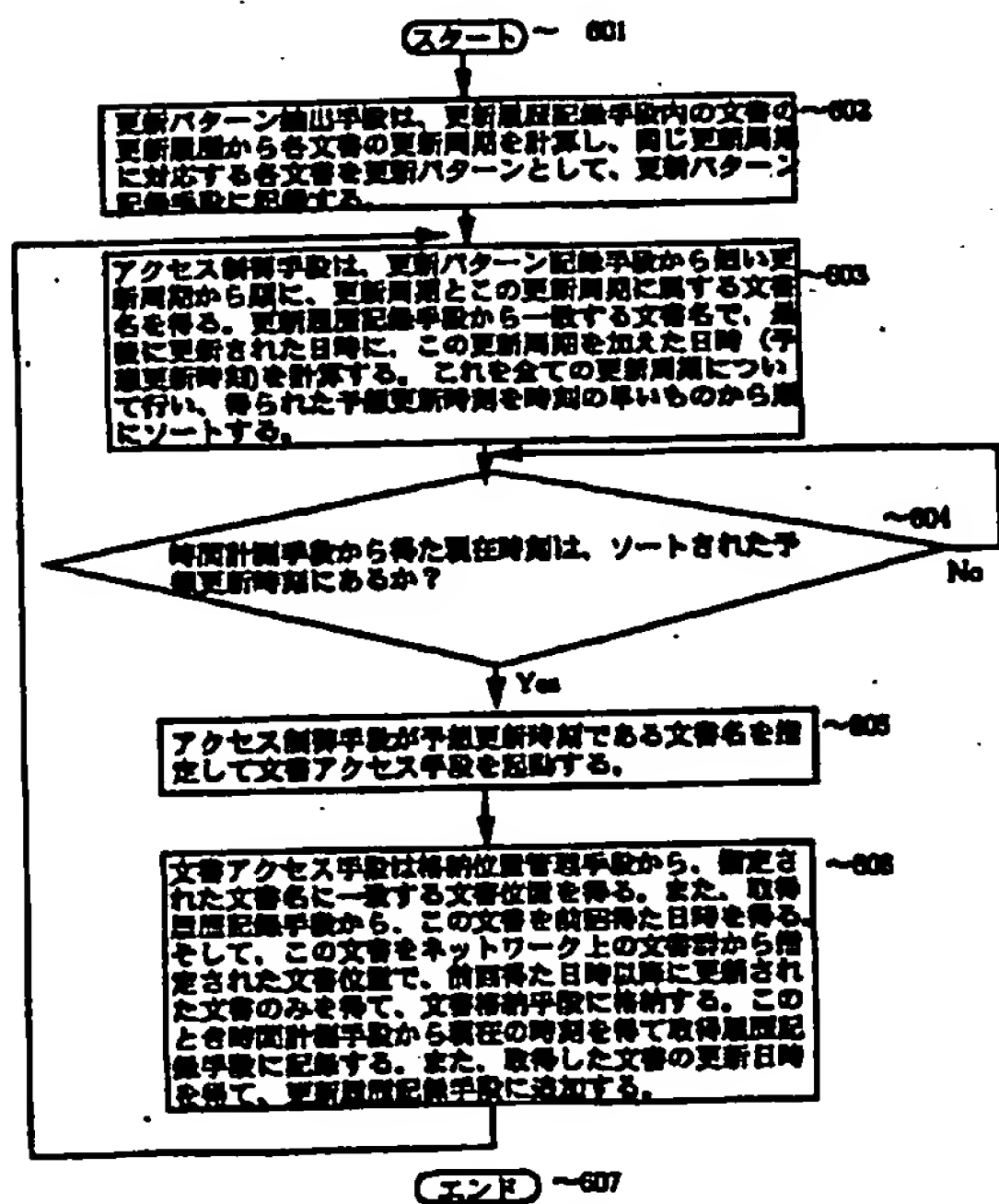




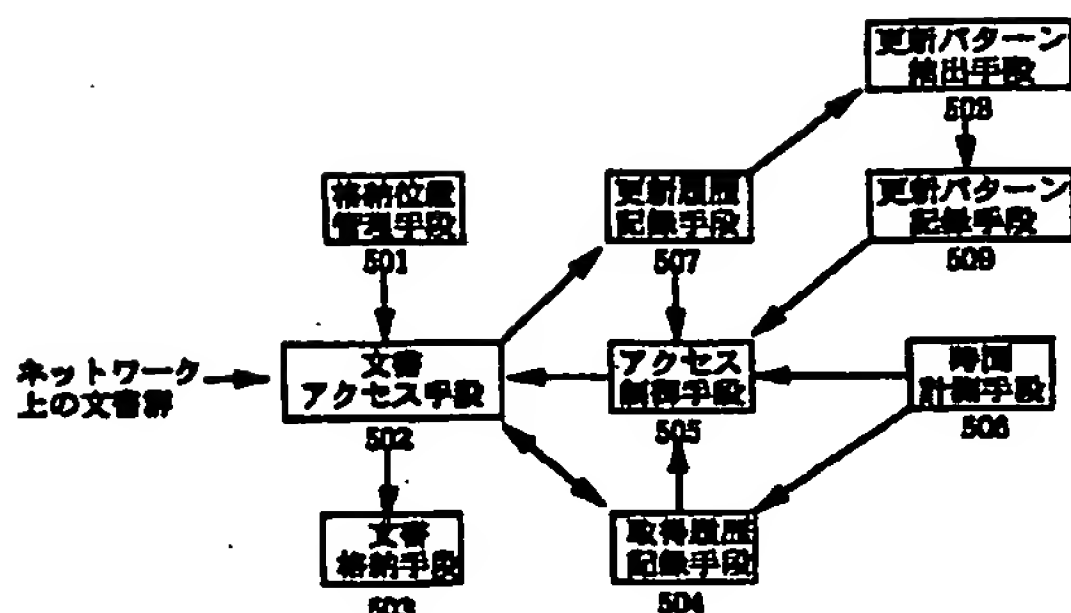
【図4】



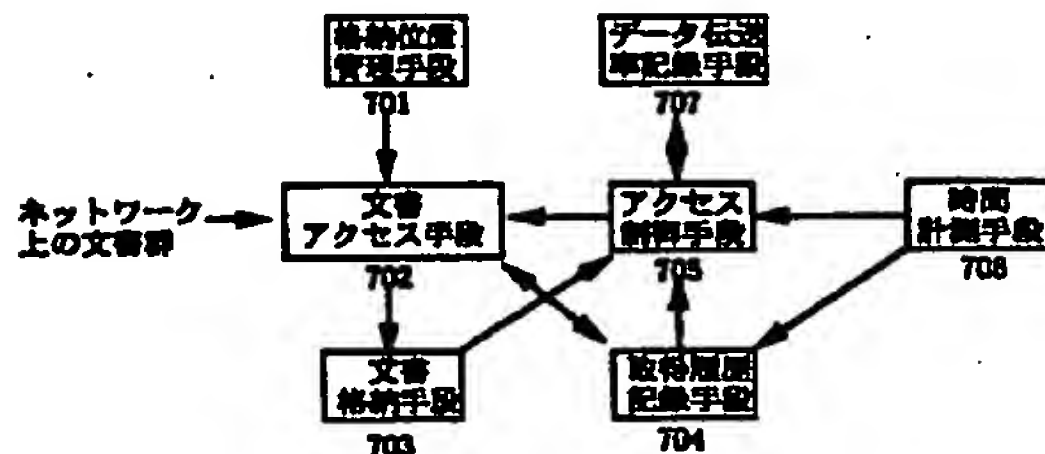
【図6】



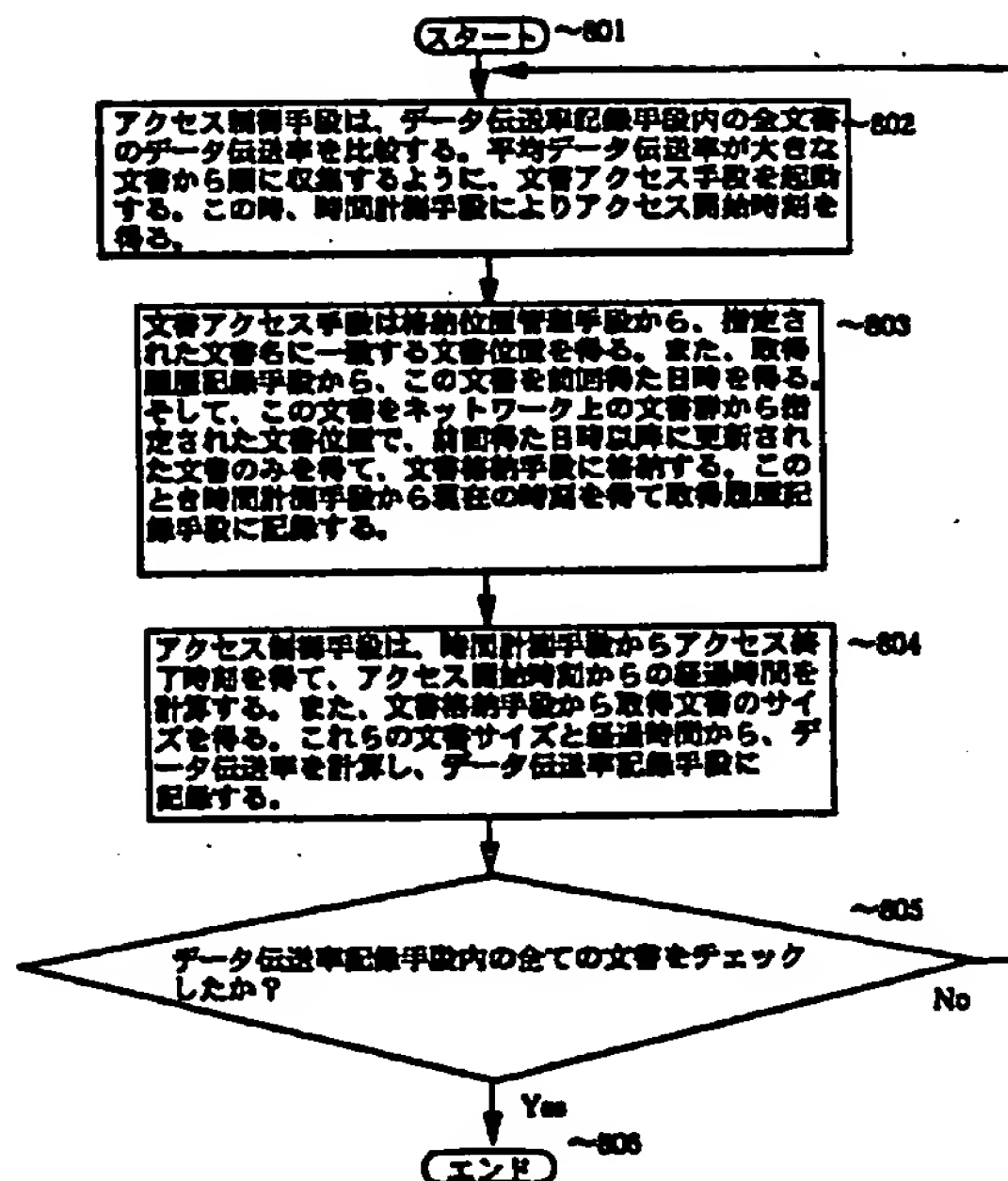
【図5】



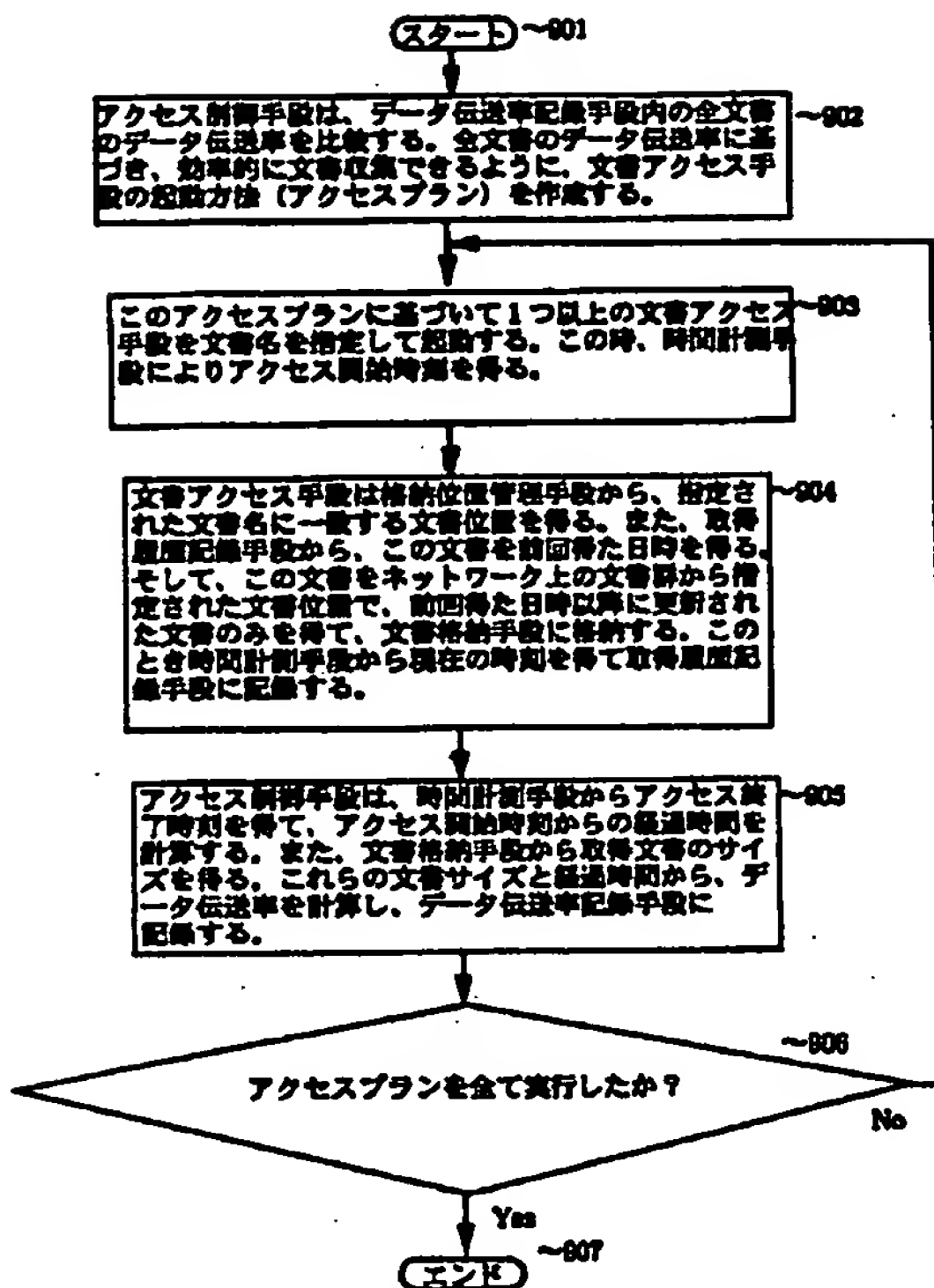
【図7】



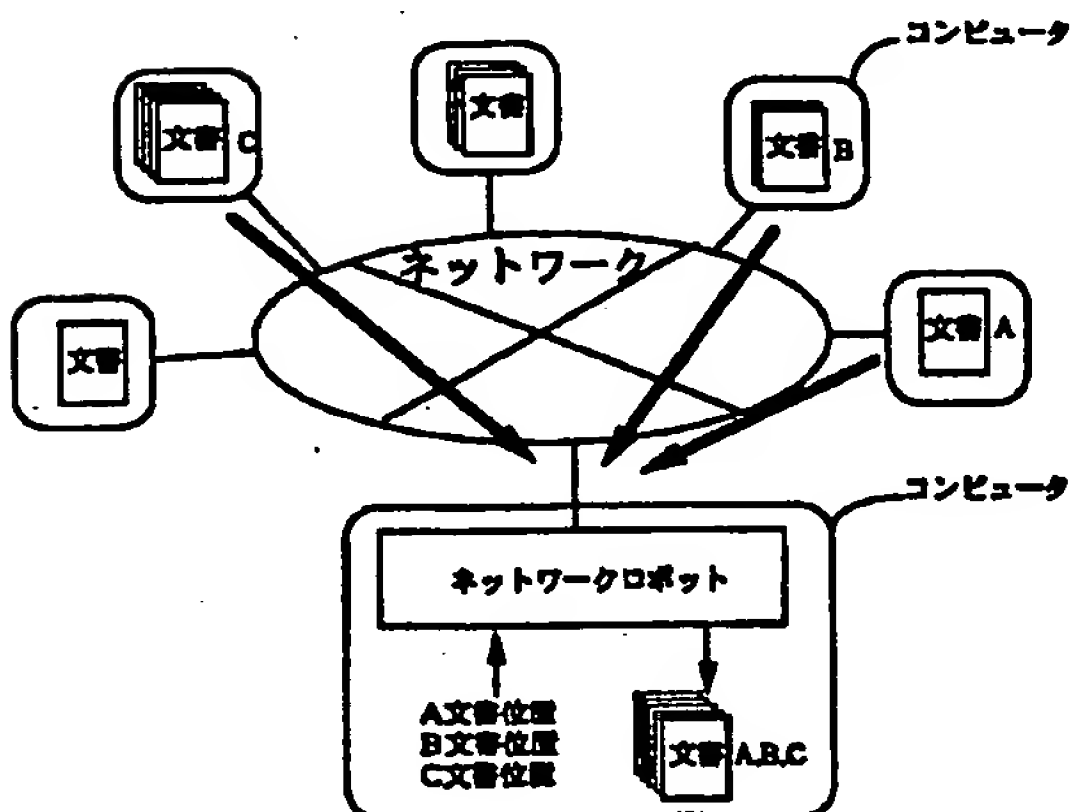
【図8】



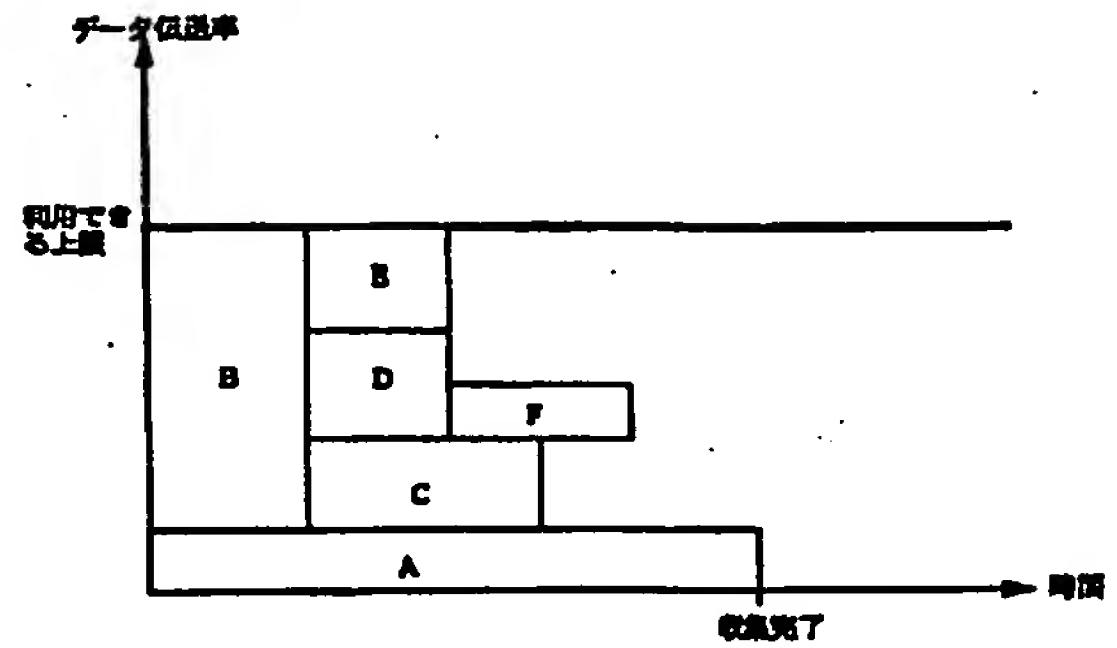
【図9】



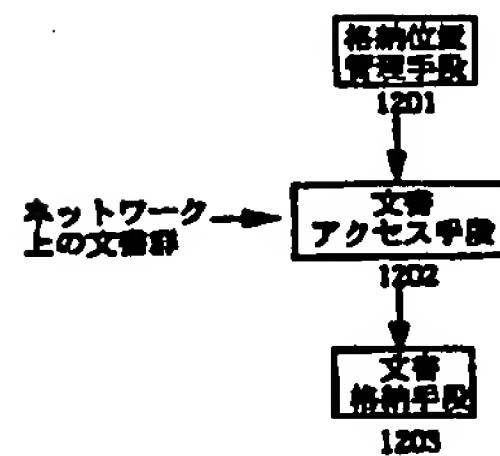
【図11】



【図10】



【図12】



フロントページの続き

(72)発明者 石川 幹人

大阪府門真市大字門真1006番地 松下電器  
産業株式会社内